

4 6. The apparatus of claim 5, wherein the network crawler is a World Wide Web robot,
5 wherein the network crawler traverses a hypertext structure of the network and retrieves the
6 content and recursively retrieves additional content referenced in the retrieved content.

7 7. The apparatus of claim 4, wherein the one or more processors, comprises:
8 a crawler processor coupled to the network crawler, wherein the crawler processor
9 receives crawling schedule information and content search criteria;

0 a network resource processor coupled to the network crawler, wherein the network
1 resource processor aggregates resource addresses of resources coupled to the one or more
2 communications networks;

3 a crawling criteria processor that compiles data related to searches to be conducted by
4 the network crawler and generates specific crawling criteria; and

5 a crawler content provider processor coupled to the network crawler that identifies,
6 tracks, indexes and ranks providers of the content, and generates content provider data,
7 wherein the network crawler receives the content provider data, the specific crawling criteria
8 and the resource addresses and crawls the network based on the received content provider
9 data, the specific crawling criteria, and the resource addresses.

8. The apparatus of claim 7, further comprising a content crawler results processor that receives content data from the network crawler, and that processes the content data and routes sorted and formatted crawling results for storage.

9. An apparatus for finding digital content in one or more communications networks,
comprising:

means for building and maintaining network resource data, wherein the network
resource data contains address data for content servers coupled to the one or more
communications networks;

means, coupled to the means for building and maintaining network resource data, for
storing the network resource data;

means for building and maintaining crawling criteria, wherein the crawling criteria
are used during a crawling operation to search for the digital content;

means for building and maintaining content provider data, wherein the content provider data comprises data related to potential providers of content on the one or more communications networks; and

means, coupled to the means for building and maintaining network resource data, the means for building and maintaining crawling criteria, and the means for building and maintaining content provider data, for crawling the communications network.

10. The apparatus of claim 9, wherein the means for building and maintaining network resource data includes means for indexing address types.

11. The apparatus of claim 10, wherein the address types include top-level domain and subdomain names, Universal Resource Identifiers, Universal Resource Locators (URLs), and Internet Protocol (IP) address numbers.

12. The apparatus of claim 10, wherein the means for indexing address types is scalable to accommodate future naming conventions.

13. The apparatus of claim 9, wherein the means for building and maintaining the network resource data includes means for updating the address data.

14. The apparatus of claim 13, wherein the means for updating the address data, comprises:

means for receiving hyperlinked domain names;

means for downloading domain name records from public and private domain name registration databases;

means for synchronizing a local Domain Name Service (DNS) database with one or more DNS databases over the one or more communications networks;

means for performing reverse domain resolution by locating URLs associated with allowable IP addressing numbers; and

means for verifying Domain Name Service aliases and duplicate URLs against IP addresses to eliminate redundant domain names.

15. The apparatus of claim 9, wherein the network resource data comprises:

URL owner identity;

URL owner contact information;

available content types;

expiration time of the domain name; and

23. The apparatus of claim 21, wherein a ranking of a content provider determines how frequently the content provider is crawled.

24. The apparatus of claim 9, wherein the means for crawling the communications network comprises one or more crawling servers, wherein the means for building and maintaining the network resource data comprises means for analyzing and subdividing the network resource data and means for providing the subdivided network resource data to the one or more crawling servers.

25. The apparatus of claim 24, wherein a crawling server comprises:

means for reading the subdivided network resource data;

means for communicating with a network resource; and

means for requesting and downloading data from the network resource.

26. The apparatus of claim 25, wherein the crawling server, further comprises:

means for comparing the content to the crawling criteria, wherein the crawling server provides data related to the content when the means for comparing indicates the crawling criteria are satisfied; and

means for following links from a first network resource to subsequent network resources, wherein the means for following links comprises:

means for analyzing hypertext structure of the first network resource to determine if the links have been crawled,

means for determining if a network resource has been downloaded or updated since a previous crawl of the network resource, and

means for analyzing the hypertext structure to determine if the link points to a network resource comprising a web page or other hypertext files.

27. The apparatus of claim 26, wherein the crawling server, further comprises:

means for caching hypertext files containing the data related to the content;

means for caching the links from the first network resource to subsequent network resources; and

means for indexing web pages and other hypertext files of interest.

28. The apparatus of claim 26, wherein the means for comparing the content to the crawling criteria comprises a comparison algorithm that compares elements in a hypertext file to the crawling criteria.

- 1 verifying DNS aliases and duplicate URLs against IP addresses; and
- 2 eliminating any duplicate URLs identified by the verifying step.
- 3 37. The method of claim 30, wherein the network resource data comprises:
- 4 URL owner identity;
- 5 URL owner contact information;
- 6 available content types;
- 7 expiration time of the domain name; and
- 8 subdomain names to be excluded during crawling.
- 9 38. The apparatus of claim 30, wherein the crawling criteria, comprises:
- 10 terms, phrases and keywords;
- 11 data type descriptions;
- 12 metadata field names; and
- 13 metadata type descriptors, wherein the metadata type descriptors are associated with
- 14 eligible content as one or more of hypertext descriptions and embedded file and data stream
- 15 attributes and metadata.
- 16 39. The apparatus of claim 30, wherein acquiring the crawling criteria comprises
- 17 automatically acquiring the crawling criteria.
- 18 40. The method of claim 39, wherein automatically acquiring the crawling criteria,
- 19 comprises:
- 20 analyzing and importing metadata schemes for standardized and proprietary content
- 21 formats;
- 22 parsing metadata field names and descriptive terms;
- 23 analyzing hypertext associated with desired hyperlinks;
- 24 analyzing text proximate to the desired hyperlinks, wherein analyzing hypertext
- 25 identifies terms that relate to a data type or content category.
- 26 41. The method of claim 30, wherein acquiring the crawling criteria comprises acquiring
- 27 the crawling criteria through manual input.
- 28 42. The method of claim 30, wherein acquiring the content provider data comprises
- 29 ranking content providers.
- 30 43. The method of claim 42, wherein a ranking of a content provider is based on one or
- 31 more of quantity of available content, provider professional association membership, amount

1 of content requested and downloaded by users of the communications network, and content
2 provider ratings, wherein the content provider ratings are provided by the users of the
3 communications network.

4 44. The method of claim 43, further comprising determining a frequency of crawling a
5 content provider based on the ranking of the content provider.

6 45. The method of claim 30, wherein crawling the network resources comprises crawling
7 with one or more crawling servers.

8 46. The method of claim 45, further comprising
9 subdividing the network resources;
10 assigning the subdivided network resources to the one or more crawling servers; and
11 at a crawler server:

12 reading data from the assigned network resources,
13 communicating with the assigned network resources,
14 downloading data from the assigned network resources.

15 47. The method of claim 46, further comprising:
16 comparing digital content from one or more of the assigned network resources to the
17 crawling criteria; and
18 acquiring data related to content that satisfies the crawling criteria.

19 48. The method of claim 46, further comprising:
20 following links from a first network resource to subsequent network resources,
21 wherein following the links comprises:

22 analyzing hypertext structure of the first network resource to determine if the
23 links have been crawled,

24 determining if a network resource has been downloaded or updated since a
25 previous crawl of the network resource, and

26 analyzing the hypertext structure to determine if the link points to a
27 network resource comprising a web page or other hypertext file.

28 49. The method of claim 48, further comprising:
29 caching hypertext files containing the data related to the content;
30 caching the links from the first network resource to subsequent network resources;

31 and

indexing web pages or other hypertext files of interest.

50. The method of claim 48, wherein comparing the content to the crawling criteria comprises using a comparison algorithm that compares elements in a hypertext file to the crawling criteria.

51. The method of claim 30, further comprising:

acquiring and processing metadata related to a network resource; and
processing content results from the crawled network resources.

52. An apparatus for controlling a remote content crawler having one or more crawling servers, the remote content crawler capable of searching one or more communications networks for data related to content available on the one or more communications networks, the apparatus, comprising:

means for communicating with components of the one or more communications networks;

means, coupled to the communications means, for executing crawling of the one or more communications networks by the remote content crawler;

means, coupled to the means for executing crawling, for routing data received by the remote content crawler; and

means, coupled to the data routing means, for aggregating data related to resources of the one or more communications networks, wherein the remote content crawler uses the aggregated data to search the one or more communications networks.

53. The apparatus of claim 52, further comprising:

means, coupled to the communications means, for building a crawling criteria database, wherein the crawling criteria comprises one or more of hypertext search guidelines, data type list, metadata search criteria, and keyword lists.

54. The apparatus of claim 52, further comprising:

means for building a content provider database, wherein data related to content providers is tracked, indexed, and ranked.

55. The apparatus of claim 52, further comprising:

means for retrieving and routing metadata related to the content available on the one or more communications networks; and

1 means, coupled to the means for retrieving and routing the metadata related to the
2 content available on the one or more communications networks, for indexing and formatting
3 the retrieved metadata.

4 56. The apparatus of claim 52, wherein the means for executing crawling, comprises:
5 means for storing data related to crawling the one or more communications networks;
6 means for initiating crawling of the one or more communications networks, the
7 means for initiating crawling comprising means for receiving administrative data related to
8 the crawling of the one or more communications networks; and
9 means for analyzing a resource data set of the one or more communications networks
10 to subdivide the resource dataset into one or more smaller resource data sets, wherein the
11 subdivision is based on one or more of overall size of the resource data set, and a number of
12 available crawling servers.

13 57. The apparatus of claim 57, wherein the means for executing crawling further
14 comprises:
15 means for determining if contents of a hypertext files meet conditions of crawling
16 criteria, comprising:
17 means for parsing the contents of the hypertext files, and
18 means for comparing the parsed content to the criteria in a criteria database,
19 wherein if a hypertext file contains sufficient matching data, the hypertext file is cached.